# Clustering in Machine Learning: Discovering Hidden Patterns in Data

Ahmed BADI

ahmedbadi905@gmail.com

linkedin.com/in/badi-ahmed

December 17, 2025

### Abstract

Clustering is a fundamental unsupervised machine learning technique that groups similar data points together without prior labeling. This article provides a comprehensive exploration of clustering methods, covering the intuition behind grouping data, mathematical foundations, and detailed explanations of major clustering algorithms including K-Means, Hierarchical Clustering, DBSCAN, and Gaussian Mixture Models. We examine the strengths and limitations of each approach, discuss distance metrics and evaluation techniques, and present real-world applications across various domains. Through clear explanations, mathematical rigor, and visual examples, this guide serves both newcomers seeking to understand clustering fundamentals and practitioners looking to apply these techniques effectively.

**Keywords:** Clustering, Unsupervised Learning, K-Means, Hierarchical Clustering, DBSCAN, Gaussian Mixture Models, Distance Metrics, Silhouette Score, Machine Learning.

## 1 Introduction

Imagine walking into a library with thousands of books scattered randomly on the floor. Your task is to organize them, but nobody tells you how. You might naturally group them by topic—fiction here, science there, history in another pile. This intuitive process of finding natural groupings is exactly what clustering does with data.

Clustering is an unsupervised machine learning technique that discovers hidden patterns and structures in unlabeled data [1]. Unlike supervised learning where we have labeled examples to learn from, clustering works with raw data and finds its own organization based on similarity.

Why is clustering so important? First, most real-world data is unlabeled—labeling is expensive and time-consuming. Second, clustering reveals insights we might not expect, discovering customer segments we didn't know existed or identifying unusual patterns in network traffic. Third, it's often a crucial first step in understanding complex datasets before applying more sophisticated techniques [2].

Clustering has transformed countless industries. Retailers use it to segment customers and personalize marketing. Biologists cluster genes to understand biological processes. Search engines cluster web pages to organize information. Social networks identify communities. Astronomers discover new types of celestial objects. The applications are endless [3].

The beauty of clustering lies in its flexibility. Different algorithms make different assumptions about data structure, allowing us to match the method to our problem. Some find spherical clusters, others discover arbitrary shapes. Some handle noise gracefully, others require clean data. Understanding these differences is key to successful application.

In this article, we'll explore the world of clustering from intuition to implementation. We'll understand what makes a good clustering, learn the mathematics behind major algorithms, see

how to evaluate results, and discover practical applications. Whether you're analyzing customer behavior, organizing documents, or exploring scientific data, this guide will equip you with both theoretical understanding and practical knowledge.

# 2 What is Clustering?

## 2.1 The Core Concept

Clustering is the task of grouping a set of objects such that objects in the same group (called a cluster) are more similar to each other than to those in other groups [1]. This similarity is typically measured using distance metrics—the closer two points are, the more similar they are.

**Formal Definition:** Given a dataset $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_n\}$ where $\mathbf{x}_i \in \mathbb{R}^d$, clustering aims to partition the data into $K$ groups $C_1, C_2, \ldots, C_K$ such that:

- Each data point belongs to at least one cluster

- Points within a cluster are similar

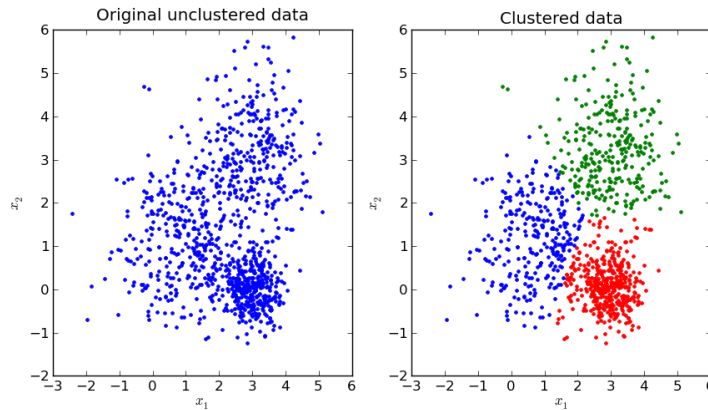- Points in different clusters are dissimilar



Figure 1: Clustering transforms unorganized data into meaningful groups.

## 2.2 Hard vs Soft Clustering

**Hard Clustering:** Each data point belongs to exactly one cluster. This is like assigning each book to a single category.

$$\mathbf{x}_i \in C_j \text{ for exactly one } j \tag{1}$$

**Soft Clustering:** Each data point has a probability or degree of membership to multiple clusters. A book might be 70% science fiction and 30% adventure.

$$\sum_{j=1}^{K} p_{ij} = 1, \quad \text{where } p_{ij} = P(\mathbf{x}_i \in C_j) \tag{2}$$

| Property | Hard Clustering | Soft Clustering |
|---|---|---|
| Membership | Binary (0 or 1) | Probabilistic (0 to 1) |
| Interpretation | Clear assignment | Captures uncertainty |
| Example | K-Means | Fuzzy C-Means, GMM |
| Use case | Clear boundaries | Overlapping clusters |

Table 1: Comparison of hard and soft clustering

## 2.3 Why Clustering Matters

Clustering addresses several fundamental challenges in data analysis:

- **Data exploration**: Understand structure in new datasets

- **Dimensionality reduction**: Represent data by cluster centers

- **Preprocessing**: Prepare data for supervised learning

- **Anomaly detection**: Identify unusual points that don't fit any cluster

- **Data compression**: Store cluster representatives instead of all points

# 3 Distance and Similarity Metrics

Before we cluster, we must define what "similar" means. This is captured by distance (or similarity) metrics.

## 3.1 Common Distance Metrics

**1. Euclidean Distance** (Most common):

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^{d}(x_i - y_i)^2} = \|\mathbf{x} - \mathbf{y}\|_2 \tag{3}$$

Measures straight-line distance. Sensitive to scale—features with larger ranges dominate.

**2. Manhattan Distance** (City Block):

$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{d}|x_i - y_i| \tag{4}$$

Sum of absolute differences. More robust to outliers than Euclidean distance.

**3. Cosine Similarity**:

$$\text{similarity}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\|\|\mathbf{y}\|} = \cos(\theta) \tag{5}$$

Measures angle between vectors, not magnitude. Perfect for text data where document length varies.

**4. Mahalanobis Distance**:

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T \Sigma^{-1} (\mathbf{x} - \mathbf{y})} \tag{6}$$

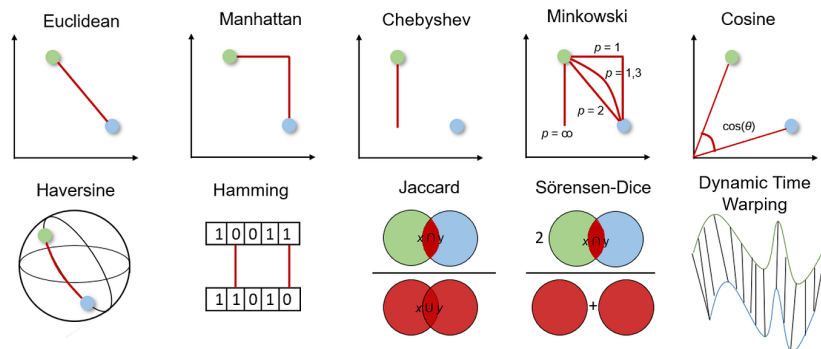Accounts for correlations between features through covariance matrix $\Sigma$.



Figure 2: Visualization of different distance metrics. Source: `https://miro.medium.com/v2/resize:fit:1400/1*BKcnB65yMzjbRAy7FQwn3w.png`

# 4 K-Means Clustering

K-Means is the most popular clustering algorithm due to its simplicity and efficiency [4].

## 4.1 The Algorithm

K-Means partitions data into $K$ clusters by minimizing the within-cluster sum of squares (WCSS):

$$\min_{\{C_k\}} \sum_{k=1}^{K} \sum_{\mathbf{x}_i \in C_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 \tag{7}$$

where $\boldsymbol{\mu}_k$ is the centroid (mean) of cluster $C_k$.
**Algorithm steps:**

1. **Initialize**: Randomly select $K$ points as initial centroids

2. **Assignment**: Assign each point to the nearest centroid

$$C_k = \{\mathbf{x}_i : \|\mathbf{x}_i - \boldsymbol{\mu}_k\| \leq \|\mathbf{x}_i - \boldsymbol{\mu}_j\| \text{ for all } j\} \tag{8}$$

3. **Update**: Recalculate centroids as the mean of assigned points

$$\boldsymbol{\mu}_k = \frac{1}{|C_k|} \sum_{\mathbf{x}_i \in C_k} \mathbf{x}_i \tag{9}$$

4. **Repeat**: Steps 2-3 until convergence (centroids don't change)

## 4.2 Choosing K: The Elbow Method

How do we choose the right number of clusters? The Elbow method plots WCSS vs $K$:

$$\text{WCSS}(K) = \sum_{k=1}^{K} \sum_{\mathbf{x}_i \in C_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 \tag{10}$$
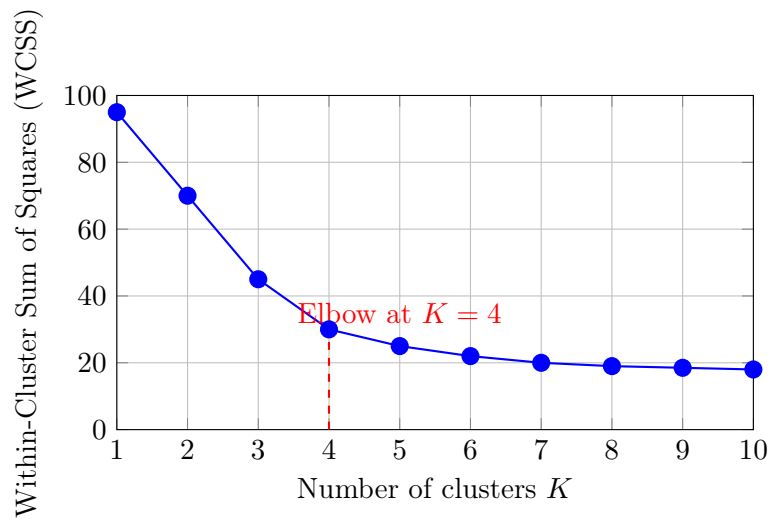


Figure 3: Elbow method: WCSS decreases rapidly until the "elbow" point, then levels off

The "elbow" indicates the optimal $K$ where adding more clusters doesn't significantly improve the model.

### 4.3    Advantages and Limitations

**Advantages:**

- Simple and intuitive

- Fast: $O(nKdi)$ where $i$ is number of iterations

- Scales well to large datasets

- Works well with spherical clusters

**Limitations:**

- Must specify $K$ in advance

- Sensitive to initialization (use K-Means++ [5])

- Assumes spherical clusters of similar size

- Sensitive to outliers

- Cannot handle non-convex shapes

# 5    Hierarchical Clustering

Hierarchical clustering builds a tree of clusters (dendrogram) without requiring $K$ upfront [6].

## 5.1    Agglomerative (Bottom-Up) Approach

Start with each point as its own cluster, then repeatedly merge the closest clusters.

**Algorithm:**

1. Start: Each point is a cluster ($n$ clusters)

2. Find the two closest clusters

3. Merge them into one cluster

4. Repeat until one cluster remains

## 5.2    Linkage Criteria

How do we measure distance between clusters? Several linkage methods exist:

**1. Single Linkage** (Minimum distance):

$$d(C_i, C_j) = \min_{\mathbf{x} \in C_i, \mathbf{y} \in C_j} d(\mathbf{x}, \mathbf{y}) \tag{11}$$

**2. Complete Linkage** (Maximum distance):

$$d(C_i, C_j) = \max_{\mathbf{x} \in C_i, \mathbf{y} \in C_j} d(\mathbf{x}, \mathbf{y}) \tag{12}$$

**3. Average Linkage**:

$$d(C_i, C_j) = \frac{1}{|C_i||C_j|} \sum_{\mathbf{x} \in C_i} \sum_{\mathbf{y} \in C_j} d(\mathbf{x}, \mathbf{y}) \tag{13}$$

**4. Ward's Method** (Minimizes variance):

$$d(C_i, C_j) = \frac{|C_i||C_j|}{|C_i| + |C_j|} \|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|^2 \tag{14}$$

## 5.3 Advantages and Limitations

**Advantages:**

- No need to specify $K$ beforehand

- Produces a hierarchy (useful for exploration)

- Deterministic (no random initialization)

- Can capture non-spherical clusters

  **Limitations:**

- Computationally expensive: $O(n^2 \log n)$ or $O(n^3)$

- Not scalable to very large datasets

- Once merged, cannot undo (greedy approach)

- Sensitive to noise and outliers

# 6 DBSCAN: Density-Based Clustering

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) discovers clusters of arbitrary shapes and identifies outliers [7].

## 6.1 Core Concepts

DBSCAN defines clusters as dense regions separated by sparse regions.

  **Key parameters:**

- $\epsilon$: Radius defining neighborhood

- MinPts: Minimum points to form a dense region

  **Point types:**

- **Core point**: Has at least MinPts within $\epsilon$

- **Border point**: Within $\epsilon$ of a core point but not core itself

- **Noise point**: Neither core nor border

## 6.2 The Algorithm

1. For each unvisited point $p$:

2. Find all points within distance $\epsilon$ (neighborhood)

3. If neighborhood has $\geq$ MinPts points, start a new cluster

4. Recursively add all density-reachable points to the cluster

5. Points not assigned to any cluster are noise

### 6.3 Advantages and Limitations

**Advantages:**

- Discovers clusters of arbitrary shapes

- Identifies outliers naturally

- No need to specify number of clusters

- Robust to noise

**Limitations:**

- Sensitive to $\epsilon$ and MinPts choice

- Struggles with varying density clusters

- High-dimensional data challenging (curse of dimensionality)

- Not deterministic on border points

# 7 Gaussian Mixture Models (GMM)

GMM is a probabilistic model assuming data is generated from a mixture of Gaussian distributions [8].

## 7.1 The Model

A GMM with $K$ components models the probability density as:

$$p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \Sigma_k) \tag{15}$$

where:

- $\pi_k$: Mixing coefficient (weight) for component $k$ with $\sum_{k=1}^{K} \pi_k = 1$

- $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \Sigma_k)$: Gaussian distribution with mean $\boldsymbol{\mu}_k$ and covariance $\Sigma_k$

## 7.2 Expectation-Maximization (EM) Algorithm

GMM parameters are estimated using the EM algorithm:

**E-step**: Calculate responsibility (posterior probability):

$$\gamma_{ik} = \frac{\pi_k \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_k, \Sigma_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_j, \Sigma_j)} \tag{16}$$

**M-step**: Update parameters:

$$N_k = \sum_{i=1}^{n} \gamma_{ik} \tag{17}$$

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{i=1}^{n} \gamma_{ik} \mathbf{x}_i \tag{18}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{i=1}^{n} \gamma_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \tag{19}$$

$$\pi_k = \frac{N_k}{n} \tag{20}$$

## 7.3 Advantages and Limitations

**Advantages:**

- Soft clustering (probabilistic memberships)

- Flexible cluster shapes (ellipsoids)

- Based on solid statistical theory

- Can incorporate prior knowledge

**Limitations:**

- Must specify number of components $K$

- Sensitive to initialization

- Computationally expensive

- Assumes Gaussian distributions

- Can converge to local optima

# 8 Evaluating Clustering Quality

How do we know if our clustering is good? Unlike supervised learning, we don't have ground truth labels.

## 8.1 Internal Evaluation Metrics

**1. Silhouette Score**:
For each point $i$, calculate:

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} \tag{21}$$

where:

- $a(i)$: Average distance to points in same cluster

- $b(i)$: Average distance to points in nearest different cluster

Range: $[-1, 1]$. Higher is better. $s(i) \approx 1$ means well-clustered.

**2. Davies-Bouldin Index**:

$$DB = \frac{1}{K} \sum_{i=1}^{K} \max_{j \neq i} \left( \frac{\sigma_i + \sigma_j}{d(\boldsymbol{\mu}_i, \boldsymbol{\mu}_j)} \right) \tag{22}$$

Lower is better. Measures average similarity between clusters.

**3. Calinski-Harabasz Index**:

$$CH = \frac{SS_B/(K-1)}{SS_W/(n-K)} \tag{23}$$

where $SS_B$ is between-cluster variance and $SS_W$ is within-cluster variance. Higher is better.

## 8.2 External Evaluation (If labels available)

**Adjusted Rand Index (ARI)**:

$$ARI = \frac{\text{RI} - E[\text{RI}]}{\max(\text{RI}) - E[\text{RI}]} \tag{24}$$

Range: $[-1, 1]$. Score of 1 means perfect clustering.

# 9 Comparison of Clustering Methods

| Algorithm | Cluster Shape | Scalability | Advantages | Need K? |
|---|---|---|---|---|
| K-Means | Spherical, equal size | Excellent | Fast, simple | Yes |
| Hierarchical | Flexible | Poor ($O(n^2)$) | No K needed, dendrogram | No |
| DBSCAN | Arbitrary | Good | Handles noise, arbitrary shapes | No |
| GMM | Ellipsoidal | Moderate | Probabilistic, flexible | Yes |

Table 2: Comparison of major clustering algorithms

# 10 Applications

## 10.1 Customer Segmentation

Companies cluster customers based on purchasing behavior, demographics, and preferences to create targeted marketing campaigns [3].

**Example features:**

- Purchase frequency

- Average transaction value

- Product categories preferred

- Age, location, income

## 10.2 Image Segmentation

Clustering pixels based on color and texture to partition images into meaningful regions for computer vision tasks.

## 10.3 Document Clustering

Grouping documents (news articles, research papers) by topic using text features like TF-IDF vectors and cosine similarity.

## 10.4 Anomaly Detection

Points that don't fit well into any cluster can be flagged as anomalies—useful for fraud detection, network intrusion detection, and quality control.

## 10.5    Recommendation Systems

Clustering users or items to find similar groups and make recommendations based on cluster behavior.

# 11    Practical Considerations

## 11.1    Data Preprocessing

**1. Feature Scaling**: Standardize features to have mean 0 and variance 1:

$$\mathbf{x}_i' = \frac{\mathbf{x}_i - \boldsymbol{\mu}}{\boldsymbol{\sigma}} \tag{25}$$

Without scaling, features with larger ranges dominate distance calculations.

**2. Dimensionality Reduction**: Use PCA or t-SNE before clustering high-dimensional data to avoid curse of dimensionality.

**3. Handle Missing Values**: Impute or remove before clustering.

## 11.2    Choosing the Right Algorithm

- **Spherical clusters, known K**: Use K-Means

- **Hierarchical structure needed**: Use Hierarchical Clustering

- **Arbitrary shapes, noisy data**: Use DBSCAN

- **Overlapping clusters, probabilistic**: Use GMM

- **Very large dataset**: Use K-Means or Mini-Batch K-Means

# 12    Conclusion

Clustering is a powerful tool for discovering hidden structure in unlabeled data. We've explored four major approaches:

- **K-Means**: Fast and simple for spherical clusters

- **Hierarchical**: Flexible with dendrogram visualization

- **DBSCAN**: Handles arbitrary shapes and noise

- **GMM**: Probabilistic with soft assignments

Each method has its strengths and trade-offs. The key to successful clustering is understanding your data, choosing appropriate distance metrics, selecting the right algorithm, and carefully evaluating results [2].

Remember: clustering is exploratory. There's often no single "correct" answer. Different clusterings can reveal different aspects of your data. Always validate results with domain knowledge and consider multiple perspectives.

As you apply clustering in practice, start simple (K-Means), visualize your data when possible, try multiple algorithms, and always preprocess appropriately. The insights you discover can transform how you understand and use your data.

# References

[1] A. K. Jain, M. N. Murty, and P. J. Flynn, "Data clustering: A review," *ACM Computing Surveys*, vol. 31, no. 3, pp. 264–323, 1999.

[2] D. Xu and Y. Tian, "A comprehensive survey of clustering algorithms," *Annals of Data Science*, vol. 2, no. 2, pp. 165–193, 2015.

[3] A. K. Jain, "Data clustering: 50 years beyond k-means," *Pattern Recognition Letters*, vol. 31, no. 8, pp. 651–666, 2010.

[4] S. P. Lloyd, "Least squares quantization in pcm," *IEEE Transactions on Information Theory*, vol. 28, no. 2, pp. 129–137, 1982.

[5] D. Arthur and S. Vassilvitskii, "K-means++: The advantages of careful seeding," in *Proceedings of the 18th Annual ACM-SIAM Symposium on Discrete Algorithms*, 2007, pp. 1027–1035.

[6] S. C. Johnson, "Hierarchical clustering schemes," *Psychometrika*, vol. 32, no. 3, pp. 241–254, 1967.

[7] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu, "A density-based algorithm for discovering clusters in large spatial databases with noise," in *Proceedings of the 2nd International Conference on Knowledge Discovery and Data Mining*, 1996, pp. 226–231.

[8] C. M. Bishop, *Pattern Recognition and Machine Learning*. New York: Springer, 2006.